

**PERBANDINGAN ALGORITMA JARO-WINKLER DISTANCE
DAN LEVENSHTEIN DISTANCE DALAM MENDETEKSI
KEMIRIPAN DOKUMEN BAHASA INDONESIA**

*Diajukan Untuk Menyusun Skripsi
di Jurusan Teknik Informatika Fakultas Ilmu Komputer UNSRI*



Oleh:

NISVA SYAKBANIA
NIM: 09021181320023

**Jurusan Teknik Informatika
FAKULTAS ILMU KOMPUTER UNIVERSITAS SRIWIJAYA
2020**

LEMBAR PENGESAHAN TUGAS AKHIR

PERBANDINGAN ALGORITMA JARO-WINKLER DISTANCE DAN LEVENSSTEIN DISTANCE DALAM MENDETEKSI KEMIRIPAN DOKUMEN BAHASA INDONESIA

Oleh :

NISVA SYAKBANIA
NIM : 09021181320023

Palembang, Agustus 2020

Pembimbing I

Pembimbing II,



Novi Yushani, M.T.
NIP. 198211082012122001



Osvari Arsalan, M.T.
NIP. 198806282018031001

Mengetahui,
Ketua Jurusan Teknik Informatika,



ii

TANDA LULUS UJIAN SIDANG SKRIPSI

Pada hari Rabu tanggal 29 Juli 2020 telah dilaksanakan ujian sidang skripsi oleh Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya.

Nama : Nisva Syakbania
NIM : 09021181320023
Judul : Perbandingan Algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance* dalam Mendeteksi Kemiripan Dokumen Bahasa Indonesia

1. Pembimbing I

Novi Yusliani, M.T.
NIP. 198211082012122001



2. Pembimbing II

Osvari Arsalan, M.T.
NIP. 198806282018031001



3. Penguji I

Samsuryadi, M.Kom., Ph.D.
NIP. 197102041997021003



4. Penguji II

Kanda Januar Miraswan, M.T.
NIP. 199001092019031012



Mengetahui,
Ketua Jurusan Teknik Informatika



HALAMAN PERNYATAAN

Yang bertanda tangan di bawah ini :

Nama : Nisva Syakbania

NIM : 09021181320023

Program Studi : Teknik Informatika

Judul Skripsi : Perbandingan Algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance* dalam Mendeteksi Kemiripan Dokumen Bahasa Indonesia

Hasil Pengecekan Software *iThenticate/Turnitin* : 2 %

Menyatakan bahwa Laporan Proyek saya merupakan hasil karya saya sendiri dan bukan hasil penjiplakan/plagiat. Apabila ditemukan unsur penjiplakan/plagiat dalam laporan proyek ini, maka saya bersedia menerima sanksi akademik dari Universitas Sriwijaya sesuai dengan ketentuan yang berlaku.

Demikian pernyataan ini saya buat dengan sebenarnya dan tidak ada paksaan oleh siapapun.

Palembang, Agustus 2020



Nisva Syakbania
NIM. 09021181320023

MOTTO DAN PERSEMBAHAN

- “Beginning always the hardest.”
- “Today, you are you that is truer than true. There is no one alive that is youer than you.” (dr. Seuss)
- “Sometimes the questions are complicated and the answers are simple.”
(dr. Seuss)
- “There is no greater weapon than a prepared mind.” (Zhuge Liang)
- “This too shall pass.”

Kupersembahkan karya tulis ini kepada :

- Myself,
- My dear parents,
- My ever-so-great sister,
- My beloved friends,
- Almamater.

**PERBANDINGAN ALGORITMA JARO-WINKLER DISTANCE
DAN LEVENSHTEIN DISTANCE DALAM MENDETEKSI
KEMIRIPAN DOKUMEN BAHASA INDONESIA**

By:
Nisva Syakbania
09021181320023

ABSTRACT

Document similarity detection is used to calculate the similarity between two or more documents based on semantic similarity or lexical similarity. This research proposed to detect similarity based on lexical similarity using a string matching techniques on each documents. Jaro-Winkler and Levenshtein Distance are algorithms usually used in string matching techniques. Jaro-Winkler Distance includes a step of calculating the length of strings in the document, counting common characters, and transposition. Levenshtein Distance is an algorithm which is used to calculate the minimum distance that needed to transform one string into the other. Testing was done with a total 19 authentic document and 6 comparative, the result of this research shows that the average error value of Levenshtein Distance is 7,86% while Jaro-Winkler Distance with average error value of 24,45%. As for computing time, four out of five testing configuration shows that Jaro-Winkler Distance have a faster computing time than Levenshtein Distance.

Keywords : Document Similarity, Jaro-Winkler Distance, Levenshtein Distance

**PERBANDINGAN ALGORITMA JARO-WINKLER DISTANCE
DAN LEVENSHTEIN DISTANCE DALAM MENDETEKSI
KEMIRIPAN DOKUMEN BAHASA INDONESIA**

Oleh:
Nisva Syakbania
09021181320023

ABSTRAK

Deteksi kemiripan dokumen digunakan untuk mengetahui kemiripan antara dua buah dokumen yang biasanya diukur berdasarkan kesamaan semantik atau kesamaan leksikal. Pada penelitian ini dilakukan deteksi kemiripan secara leksikal dengan cara melakukan pencocokan string pada setiap dokumen. Algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance* merupakan algoritma yang sering digunakan dalam teknik pencocokan pola. Algoritma *Jaro-Winkler Distance* menghitung panjang kata dalam dokumen, kata yang sama, dan jumlah transposisi. Sedangkan algoritma *Levenshtein Distance* menghitung jarak yang dibutuhkan untuk mengubah satu kata menjadi kata lain. Pengujian dilakukan dengan menggunakan 19 dokumen asli dan 6 dokumen pembanding menghasilkan rata-rata nilai error sebesar 7,86% untuk algoritma *Levenshtein Distance* dan 24,45% untuk algoritma *Jaro-Winkler Distance*. Untuk waktu proses, empat dari lima skenario pengujian menunjukkan *Jaro-Winkler Distance* memiliki waktu komputasi yang lebih cepat daripada *Levenshtein Distance*.

Kata Kunci: Kemiripan Dokumen, *Jaro-Winkler Distance*, *Levenshtein Distance*

KATA PENGANTAR

Segala puji dan syukur bagi Allah SWT atas berkah dan rahmat-Nya yang telah diberikan kepada Penulis sehingga dapat menyelesaikan Tugas Akhir ini dengan baik. Tugas akhir yang berjudul “**Perbandingan Algoritma Jaro-Winkler Distance dan Levenshtein Distance dalam Mendeteksi Kemiripan Dokumen Bahasa Indonesia**” ini disusun untuk memenuhi salah satu syarat guna menyelesaikan pendidikan program Strata-1 pada Fakultas Ilmu Komputer Program Studi Teknik Informatika di Universitas Sriwijaya.

Pada kesempatan ini Penulis ingin menyampaikan ucapan terima kasih kepada pihak-pihak yang telah memberikan dukungan, bimbingan, dan motivasi selama proses penelitian ini dilaksanakan. Secara khusus Penulis ingin menyampaikan terima kasih kepada:

1. Allah SWT yang selalu menjawab doa hamba-Nya.
2. Nisva Syakbania, terima kasih karena tidak menyerah dan selalu percaya pada diri sendiri bahwa ia bisa menyelesaikan tugas akhirnya.
3. Kedua orang tua tercinta, terima kasih untuk selalu menyebut namaku dalam setiap doa. Keyakinan dan dukungan moril maupun materil yang tiada henti telah membantu Penulis dalam menyelesaikan Tugas Akhir ini.
4. Saudaraku, Mega Puspita, atas dukungan dan kata-kata bijak yang telah membangkitkan motivasi Penulis.
5. Bapak Jaidan Jauhari, M.T. selaku Dekan Fakultas Ilmu Komputer Universitas Sriwijaya.

6. Bapak Rifkie Primartha, M.T. selaku Ketua Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya.
7. Ibu Novi Yusliani, M.T. selaku dosen pembimbing I dan Bapak Osvari Arsalan, M.T. selaku dosen pembimbing II yang telah membimbing, mengarahkan, dan memberikan motivasi untuk Penulis dalam proses penggerjaan Tugas Akhir.
8. Bapak Samsuryadi, M.Kom., Ph.D. dan Bapak Kanda Januar, M.T. selaku penguji Tugas Akhir yang telah memberi nasihat dan saran yang membangun.
9. Bapak Julian Supardi dan Bapak Danny Matthew Saputra, M.Sc. selaku pembimbing akademik.
10. Seluruh dosen Program Studi Teknik Informatika yang telah memberikan bekal ilmu selama masa perkuliahan.
11. Sahabat seperjuangan Teknik Informatika Angkatan 2013.
12. Seluruh staff administrasi dan pengawai yang telah membantu dalam urusan administrasi

Penulis menyadari masih terdapat banyak kekurangan dalam penyusunan Tugas Akhir disebabkan oleh keterbatasan pengetahuan dan pengalaman. Oleh karena itu kritik dan saran yang membangun sangat diharapkan untuk menyempurnakan Tugas Akhir ini. Semoga Tugas Akhir ini dapat bermanfaat bagi kita semua.

Palembang, Agustus 2020

Nisva Syakbania

DAFTAR ISI

Halaman

HALAMAN JUDUL.....	i
HALAMAN PENGESAHAN	ii
HALAMAN TANDA LULUS SIDANG SKRIPSI	iii
HALAMAN PERNYATAAN	iii
HALAMAN MOTTO DAN PERSEMBAHAN.....	v
ABSTRACT.....	vi
ABSTRAK	vii
KATA PENGANTAR	viii
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiv
DAFTAR GAMBAR	xviii
DAFTAR LAMPIRAN	xx
BAB I PENDAHULUAN.....	I-1
1.1 Pendahuluan	I-1
1.2 Latar Belakang Masalah.....	I-1
1.3 Rumusan Masalah	I-4
1.4 Tujuan Penelitian	I-4
1.5 Manfaat Penelitian	I-5
1.6 Batasan Masalah.....	I-5
1.7 Sistematika Penulisan.....	I-5
1.8 Kesimpulan	I-7
BAB II TINJAUAN PUSTAKA.....	II-1
2.1 Pendahuluan	II-1
2.2 Kemiripan Teks.....	II-1
2.3 Prapengolahan Teks	II-2
2.3.1 Segmentasi Kalimat	II-2
2.3.2 Case Folding	II-3

2.3.3 Tokenizing	II-4
2.3.4 Stopword Removal	II-5
2.3.5 Stemming.....	II-6
2.4 Jaro-Winkler Distance.....	II-7
2.5 Levenshtein Distance	II-10
2.6 Penelitian Terkait	II-13
2.7 Kesimpulan	II-14
 BAB III METODOLOGI PENELITIAN.....	III-1
3.1 Pendahuluan	III-1
3.2 Pengumpulan Data	III-1
3.2.1 Jenis Data.....	III-1
3.2.2 Sumber Data	III-1
3.2.3 Metode Pengumpulan Data.....	III-2
3.3 Tahapan Penelitian	III-2
3.3.1 Langkah-langkah Penelitian	III-2
3.3.2 Kerangka Kerja.....	III-4
3.3.3 Kriteria Pengujian	III-8
3.3.4 Format Data Pengujian	III-9
3.3.5 Alat yang Digunakan dalam Pelaksanaan Penelitian.....	III-10
3.3.6 Melakukan Analisis Hasil Pengujian dan Membuat Kesimpulan	III-11
3.4 Metode Pengembangan Perangkat Lunak	III-12
3.5 Managemen Proyek Penelitian.....	III-15
 BAB IV PENGEMBANGAN PERANGKAT LUNAK	IV-1
4.1 Pendahuluan	IV-1
4.2 Fase Insepsi	IV-1
4.2.1 Pemodelan Bisnis	IV-1
4.2.2 Kebutuhan Sistem	IV-2
4.2.3 Analisis dan Desain	IV-3
4.2.3.1 Analisis Kebutuhan Perangkat Lunak	IV-3
4.2.3.2 Analisis Data	IV-5
4.2.3.3 Analisis Prapengolahan	IV-8

4.2.3.4 Analisis Perhitungan Kemiripan dengan Jaro-Winkler Distance	IV-14
4.2.3.5 Analisis Perhitungan Kemiripan dengan Levenshtein Distance	IV-17
4.2.4 Desain Perangkat Lunak	IV-23
4.2.4.1 Model Use Case	IV-23
4.2.4.2 Diagram Kelas Analisis.....	IV-31
4.3 Fase Elaborasi	IV-40
4.3.1 Pemodelan Bisnis	IV-40
4.3.1.1 Perancangan Data	IV-40
4.3.1.2 Perancangan Antarmuka	IV-40
4.3.2 Kebutuhan Sistem.....	IV-42
4.3.3 Diagram Sequence	IV-43
4.4 Fase Konstruksi.....	IV-49
4.4.1 Diagram Kelas	IV-49
4.4.2 Implementasi.....	IV-51
4.4.2.1 Implementasi Kelas	IV-51
4.4.2.2 Implementasi Antarmuka	IV-54
4.5 Fase Transisi	IV-56
4.5.1 Pemodelan Bisnis	IV-56
4.5.2 Kebutuhan Sistem.....	IV-57
4.5.3 Rencana Pengujian Black Box.....	IV-57
4.5.3.1 Rencana Pengujian Use Case Memasukkan Dokumen Asli	IV-57
4.5.3.2 Recana Pengujian Use Case Memasukkan Dokumen Pembanding.....	IV-58
4.5.3.4 Rencana Pengujian Use Case Melakukan Deteksi Kemiripan dengan Jaro-Winkler Distance	IV-59
4.5.3.5 Rencana Pengujian Use Case Melakukan Deteksi Kemiripan dengan Levenshtein Distance.....	IV-59
4.5.4 Implementasi Pengujian Black Box.....	IV-60
4.5.5 Rencana Pengujian White Box	IV-66
4.5.6 Implementasi Pengujian White Box	IV-67
4.6 Kesimpulan	IV-74

BAB V HASIL DAN ANALISIS PENELITIAN.....	V-1
5.1 Pendahuluan	V-1
5.2 Data Hasil Penelitian.....	V-1
5.2.1 Data Hasil Skenario I.....	V-2
5.2.2 Data Hasil Skenario II	V-3
5.2.3 Data Hasil Skenario III	V-4
5.2.4 Data Hasil Skenario IV	V-5
5.2.5 Data Hasil Skenario V	V-6
5.3 Analisis Hasil Penelitian	V-7
5.3.1 Analisis Akurasi.....	V-7
5.3.2 Analisis Waktu Proses	V-10
5.4 Kesimpulan	V-12
 BAB VI KESIMPULAN DAN SARAN	VI-1
6.1 Pendahuluan	VI-1
6.2 Kesimpulan	VI-1
6.3 Saran.....	VI-2
 DAFTAR PUSTAKA	xxi
LAMPIRAN	xxii

DAFTAR TABEL

	Halaman
Tabel II-1. Segmentasi Kalimat	II-2
Tabel II-2. Casing.....	II-3
Tabel II-3. Tokenizing.....	II-4
Tabel II-4. Stopword Removal.....	II-5
Tabel II-5. Representasi Panjang String dalam Bentuk Matriks.....	II-11
Tabel II-6. Matriks Perhitungan Levenshtein Distance.....	II-12
Tabel III- 1. Kategori Dokumen.....	III-5
Tabel III- 2. Format Data Pengujian.....	III-9
Tabel III- 3. Spesifikasi Kebutuhan Perangkat Keras dan Perangkat Lunak....	III-10
Tabel III- 4. Rata-Rata Nilai Error.....	III-11
Tabel III- 5. Rata-Rata Waktu Proses.....	III-12
Tabel IV- 1. Kebutuhan Fungsional.....	IV-4
Tabel IV- 2. Kebutuhan Non-Fungsional.....	IV-4
Tabel IV- 3. Daftar Dokumen untuk Skenario Pengujian.....	IV-5
Tabel IV- 4. Daftar Dokumen Beda Struktur Kalimat.....	IV-7
Tabel IV- 5. Contoh Dokumen Bahasa Indonesia.....	IV-8
Tabel IV- 6. Hasil Segmentasi Kalimat dari Contoh Dokumen Bahasa Indonesia.....	IV-9
Tabel IV- 7. Hasil Casing dari Contoh Dokumen Bahasa Indonesia.....	IV-10
Tabel IV- 8. Hasil Tokenizing dari Contoh Dokumen Bahasa Indonesia.....	IV-11

Tabel IV- 9. Hasil Stopword Removal dari Contoh Dokumen Bahasa Indonesia.....	IV-12
Tabel IV- 10. Hasil Stemming dari Contoh Dokumen Bahasa Indonesia.....	IV-13
Tabel IV- 11. Hasil Perhitungan Jumlah Token yang Sama.....	IV-15
Tabel IV- 12. Matriks Hasil Perhitungan Levenshtein Distance.....	IV-21
Tabel IV- 13. Definisi Aktor Use Case.....	IV-24
Tabel IV- 14. Definisi Use Case.....	IV-24
Tabel IV- 15. Skenario Use Case Memasukkan Dokumen Asli.....	IV-26
Tabel IV- 16. Skenario Use Case Memasukkan Dokumen Pembanding.....	IV-27
Tabel IV- 17. Skenario Use Case Melakukan Prapengolahan Data.....	IV-28
Tabel IV- 18. Skenario Use Case Melakukan Deteksi Kemiripan dengan Jaro-Winkler Distance.....	IV-29
Tabel IV- 19. Skenario Use Case Melakukan Deteksi Kemiripan dengan Levenshtein Distance.....	IV-30
Tabel IV- 20. Implementasi Kelas.....	IV-51
Tabel IV- 21. Rencana Pengujian Use Case Memasukkan Dokumen Asli.....	IV-58
Tabel IV- 22. Rencana Pengujian Use Case Memasukkan Dokumen Pembanding.....	IV-58
Tabel IV- 23. Rencana Pengujian Use Case Melakukan Prapengolahan Data	IV-58
Tabel IV- 24. Rencana Pengujian Use Case Melakukan Deteksi Kemiripan dengan Jaro-Winkler Distance.....	IV-59
Tabel IV- 25. Rencana Pengujian Use Case Melakukan Deteksi Kemiripan dengan Levenshtein Distance.....	IV-59
Tabel IV- 26. Pengujian Use Case Memasukkan Dokumen Asli.....	IV-61

Tabel IV- 27. Pengujian Use Case Memasukkan Dokumen Pembanding.....	IV-62
Tabel IV- 28. Pengujian Use Case Melakukan Prapengolahan Data.....	IV-63
Tabel IV- 29. Pengujian Use Case Melakukan Deteksi Kemiripan dengan Jaro-Winkler Distance.....	IV-64
Tabel IV- 30. Pengujian Use Case Melakukan Deteksi Kemiripan dengan Levenshtein Distance.....	IV-65
Tabel IV- 31. Pengujian White Box Skenario Dokumen 100% copy&paste ..	IV-68
Tabel IV- 32. Pengujian White Box untuk Skenario Dokumen 50% copy&paste.....	IV-69
Tabel IV- 33. Pengujian White Box untuk Skenario Dokumen 20% copy&paste.....	IV-70
Tabel IV- 34. Pengujian White Box untuk Skenario Dokumen Gabungan.....	IV-71
Tabel IV- 35. Pengujian White Box untuk Skenario Dokumen Beda Struktur	IV-72
Tabel IV- 36. Pengujian White Box untuk Skenario Dokumen Beda Isi.....	IV-73
Tabel V- 1. Hasil Deteksi Kemiripan Dokumen Skenario I.....	V-2
Tabel V- 2. Hasil Deteksi Kemiripan Dokumen Skenario II.....	V-3
Tabel V- 3. Hasil Deteksi Kemiripan Dokumen Skenario III.....	V-4
Tabel V- 4. Hasil Deteksi Kemiripan Dokumen Skenario IV.....	V-5
Tabel V- 5. Hasil Deteksi Kemiripan Dokumen Skenario V.....	V-6
Tabel V- 6. Rata-Rata Nilai Error.....	V-8
Tabel V- 7. Rata-Rata Waktu Proses.....	V-10

DAFTAR GAMBAR

Gambar III-1. Diagram Blok Tahapan Penelitian.....	III-3
Gambar III- 2. Diagram Tahapan Perangkat Lunak.....	III-4
Gambar III- 3. Diagram Blok Proses Perhitungan Kemiripan dengan Jaro-Winkler Distance.....	III-7
Gambar III- 4. Diagram blok proses perhitungan kemiripan dengan Levenshtein Distance.....	III-8
Gambar III- 5. Penjadwalan untuk Tahap Menentukan Ruang Lingkup dan Unit Penelitian.....	III-15
Gambar III- 6. Penjadwalan untuk Tahap Menentukan Dasar Teori dan Menentukan Kriteria Pengujian.....	III-15
Gambar III- 7. Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Insepsi.....	III-16
Gambar III- 8. Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Elaborasi.....	III-16
Gambar III- 9. Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Konstruksi.....	III-16
Gambar III- 10. Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Transisi.....	III-17
Gambar III- 11. Penjadwalan untuk Tahap Melakukan Pengujian, Analisa Hasil Pengujian dan Membuat Kesimpulan.....	III-17
Gambar IV- 1. Contoh Dokumen Beda Struktur Kalimat.....	IV-8
Gambar IV- 2. Diagram Use Case.....	IV-23

Gambar IV- 3. Kelas Analisis Memasukkan Dokumen Asli.....	IV-31
Gambar IV- 4. Kelas Analisis Memasukkan Dokumen Pembanding.....	IV-32
Gambar IV- 5. Kelas Analisis Melakukan Prapengolahan Data.....	IV-34
Gambar IV- 6. Kelas Analisis Melakukan Deteksi Kemiripan Dokumen dengan Jaro-Winkler Distance.....	IV-35
Gambar IV- 7. Kelas Analisis Melakukan Deteksi Kemiripan Dokumen dengan Levenshtein Distance.....	IV-36
Gambar IV- 8. Diagram Aktivitas Memasukkan Dokumen Asli.....	IV-37
Gambar IV- 9. Diagram Aktivitas Memasukkan Dokumen Pembanding.....	IV-37
Gambar IV- 10. Diagram Aktivitas Melakukan Prapengolahan Data.....	IV-38
Gambar IV- 11. Diagram Aktivitas Melakukan Deteksi Kemiripan dengan Jaro-Winkler Distance.....	IV-39
Gambar IV- 12. Diagram Aktivitas Melakukan Deteksi Kemiripan dengan Levenshtein Distance.....	IV-39
Gambar IV- 13. Perancangan Antarmuka Jaro-Winkler Distance.....	IV-41
Gambar IV- 14. Perancangan Antarmuka Levenshtein Distance.....	IV-41
Gambar IV- 15. Diagram Sequence Memasukkan Dokumen Asli.....	IV-44
Gambar IV- 16. Diagram Sequence Memasukkan Dokumen Pembanding.....	IV-45
Gambar IV- 17. Diagram Sequence Melakukan Prapengolahan Data.....	IV-46
Gambar IV- 18. Diagram Sequence Melakukan Deteksi Kemiripan dengan Jaro-Winkler Distance.....	IV-47
Gambar IV- 19. Diagram Sequence Melakukan Deteksi Kemiripan dengan Levenshtein Distance.....	IV-48
Gambar IV- 20. Diagram Kelas Perangkat Lunak.....	IV-50

Gambar IV- 21. Antarmuka MainMenu.....	IV-55
Gambar IV- 22. Antarmuka JaroWinklerWindow.....	IV-55
Gambar IV- 23. Antarmuka LevenshteinWindow.....	IV-56
Gambar IV- 24. Diagram Alur Pengujian White Box Algoritma Jaro-Winkler Distance.....	IV-66
Gambar IV- 25. Diagram Alur Pengujian White Box Algoritma Levenshtein Distance.....	IV-67
Gambar IV- 26. Hasil Pengujian White Box Skenario Dokumen 100% copy&paste.....	IV-68
Gambar IV- 27. Hasil Pengujian White Box Skenario Dokumen 50% copy&paste.....	IV-69
Gambar IV- 28. Hasil Pengujian White Box Skenario Dokumen 20% copy&paste.....	IV-70
Gambar IV- 29. Hasil Pengujian White Box Skenario Dokumen Gabungan..	IV-71
Gambar IV- 30. Hasil Pengujian White Box Skenario Dokumen Beda Struktur.....	IV-72
Gambar IV- 31. Hasil Pengujian White Box Skenario Dokumen Beda I.....	IV-73
Gambar V- 1. Grafik Rata-Rata Nilai Error.....	V-8
Gambar V- 2. Grafik Rata-Rata Waktu Proses.....	V-11

DAFTAR LAMPIRAN

Halaman

Lampiran I: Perhitungan Manual Skenario Pengujian.....	L-1
A. Perhitungan Manual Skenario Pengujian I (100% copy&paste).....	L-1
B. Perhitungan Manual Skenario Pengujian III (50% copy&paste).....	L-3
C. Perhitungan Manual Skenario Pengujian III (20% copy&paste).....	L-5
D. Perhitungan Manual Skenario Pengujian IV (Gabungan).....	L-7
E. Perhitungan Manual Skenario Pengujian V (Beda Struktur).....	L-10
F. Perhitungan Manual Skenario Pengujian VI (Beda Kalimat).....	L-12
Lampiran II: DATASET.....	L-15
Lampiran III: Kode Program.....	L-38

BAB I

PENDAHULUAN

1.1 Pendahuluan

Pada bab ini membahas secara sistematis mengenai penelitian tugas akhir yang meliputi latar belakang masalah, rumusan masalah, tujuan penelitian, manfaat penelitian, dan batasan masalah yang menjelaskan batasan perangkat lunak yang akan dibuat.

1.2 Latar Belakang Masalah

Perkembangan pesat di dunia teknologi dan informasi memudahkan setiap kalangan untuk mengakses berbagai informasi. Seiring dengan kemudahan akses informasi ini terdapat juga beberapa dampak negatif, salah satunya praktik pelanggaran etika akademis yang dilakukan oleh kalangan akademisi. Para akademisi ini dapat dengan mudah mengambil data dan informasi yang tersebar di internet lalu menggunakan tanpa mencantumkan sumber. Hal ini dibuktikan dengan banyak ditemukannya kasus plagiarisme oleh mahasiswa dan tenaga pengajar di universitas.

Kemiripan teks merupakan salah satu bentuk plagiarisme dimana aspek yang ditiru berupa kata demi kata, kalimat, paragraf atau bahkan keseluruhan karya orang lain (Sastroasmoro, 2007). Pendekripsiannya kemiripan ini dapat dilakukan dengan melakukan perbandingan antara teks asli dengan teks yang dicurigai.

Apabila tingkat kemiripan teks mencapai persentase tertentu (tergantung kebijakan tiap instansi), maka karya tulis tersebut dapat dinyatakan sebagai plagiarisme.

Salah satu langkah yang dapat dilakukan untuk mendeteksi kemiripan teks adalah dengan deteksi secara manual dan otomatis. Deteksi kemiripan dokumen secara otomatis sendiri sudah banyak dikembangkan dalam beberapa tahun ini. Deteksi kemiripan teks dengan membandingkan keseluruhan isi dokumen kata per kata dapat dilakukan dengan pendekatan pencocokan string, salah satunya *approximate string matching*. Metode ini merupakan teknik dalam pencocokan pola pada string berdasarkan kemiripan tekstual atau penulisan meliputi jumlah karakter dan susunan karakter dalam dokumen (Rochmawati & Kusumaningrum, 2016). Contoh algoritma yang sering digunakan dalam pendekatan ini adalah algoritma *brute force*, *boyer-moore*, *jaro-winkler*, *levenshtein distance*, dan *hamming distance*.

Pendeteksian kemiripan dokumen bahasa Indonesia pada penelitian ini akan menggunakan algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance*. Algoritma *Jaro-Winkler Distance* adalah algoritma untuk menghitung kesamaan antara dua string dan merupakan salah satu algoritma yang paling efektif dalam menghitung tingkat kemiripan. Dalam penelitian perbandingan metode *Approximate String Matching* dalam mengidentifikasi kesalahan penulisan (Rochmawati & Kusumaningrum, 2016), algoritma *Jaro-Winkler Distance* merupakan algoritma yang memiliki nilai terbaik dalam melakukan perbaikan ejaan pada bahasa Indonesia dengan nilai MAP (*Mean Average Precision*) sebesar 0,87 dibandingkan algoritma lainnya. Dalam penelitian yang dilakukan oleh (Tinaliah &

Elizabeth, 2018) yang berjudul “*Perbandingan Hasil Deteksi Plagiarisme Dokumen dengan Metode Jaro-Winkler Distance dan Metode Latent Semantic Analysis*”, metode Jaro-Winkler memberikan hasil plagiat 100% pada data yang sama persis dibandingkan dengan metode *LSA* yang menghasilkan nilai kemiripan 97,14%.

Levenshtein Distance merupakan algoritma pengukuran string untuk menghitung perbedaan antara dua kata dimana semakin kecil nilai yang dimiliki antara dua kata tersebut maka semakin tinggi tingkat kemiripan. Algoritma *Levenshtein Distance* sudah digunakan di berbagai bidang seperti perbaikan ejaan, mesin pencari, dan juga deteksi plagiarisme. Berdasarkan hasil penelitian terdahulu, didapatkan hasil akurasi sebesar 86% dalam pengoreksian kesalahan masukan pada sistem pencarian (Yulianto, Arifudin, & Alamsyah, 2018).

1.3 Rumusan Masalah

Berdasarkan latar belakang yang telah dikemukakan sebelumnya, maka penelitian tugas akhir yang dilakukan adalah mendeteksi kemiripan dokumen Bahasa Indonesia menggunakan algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance*. Untuk mendukung rumusan masalah tersebut, maka penelitian ini akan dibagi dalam *research questions*, yaitu:

1. Berapa lama waktu proses deteksi kemiripan dokumen bahasa Indonesia dengan algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance*.

2. Berapa nilai kemiripan yang dihasilkan algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance* dalam mendeteksi kemiripan dokumen bahasa Indonesia.

1.4 Tujuan Penelitian

Tujuan penelitian ini adalah:

1. Untuk mengetahui waktu proses algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance* dalam mendeteksi kemiripan dokumen bahasa Indonesia.
2. Untuk mengetahui perbedaan persentase kemiripan dokumen yang dihasilkan dari algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance*.

1.5 Manfaat Penelitian

Manfaat penelitian ini adalah:

1. Hasil penelitian ini dapat mendeteksi kemiripan antara dua dokumen yang kemudian dapat digunakan sebagai dugaan awal dalam deteksi plagiat.
2. Hasil penelitian ini dapat digunakan sebagai referensi untuk pemilihan algoritma yang lebih baik untuk mendeteksi kemiripan dokumen.

1.6 Batasan Masalah

Batasan masalah pada penelitian ini adalah:

1. Dataset yang digunakan merupakan dokumen dalam bahasa Indonesia.
2. Format dokumen masukan adalah *.txt dan *.pdf.

3. Penelitian ini hanya mengukur persentase kemiripan antar dua buah dokumen.

1.7 Sistematika Penulisan

Sistematika penulisan skripsi ini adalah sebagai berikut:

BAB I. PENDAHULUAN

Pada bab ini diuraikan mengenai latar belakang, perumusan masalah, tujuan dan manfaat penelitian, batasan masalah, metodologi penelitian, dan sistematika penulisan.

BAB II. KAJIAN LITERATUR

Pada bab ini akan membahas landasan teori yang digunakan dalam penelitian ini, seperti definisi kemiripan teks, analisis *preprocessing*, analisis algoritma *Jaro-Winkler Distance* dan algoritma *Levenshtein Distance*. Selain itu akan dibahas mengenai penelitian-penelitian terkait yang penelitian ini.

BAB III. METODOLOGI PENELITIAN

Pada bab ini akan dibahas mengenai tahapan yang akan dilaksanakan pada penelitian ini. Masing-masing rencana tahapan penelitian dideskripsikan dengan rinci berdasarkan pada kerangka kerja. Di akhir bab akan dijabarkan perancangan manajemen proyek perangkat lunak untuk pelaksanaan penelitian ini.

BAB IV. PENGEMBANGAN PERANGKAT LUNAK

Pada bab ini akan dibahas mengenai analisis, perancangan dan lingkungan implementasi perangkat lunak deteksi kemiripan dokumen bahasa Indonesia, implementasi algoritma *Jaro-Winkler Distance*, implementasi algoritma *Levenshtein Distance*, hasil eksekusi, dan hasil pengujian.

BAB V. HASIL DAN ANALISIS PENELITIAN

Pada bab ini pengujian dilakukan berdasarkan skenario-skenario pengujian yang telah dirancang sebelumnya. Kemudian dilakukan analisis sebagai basis dari kesimpulan yang diambil dalam penelitian ini.

BAB VI. KESIMPULAN DAN SARAN

Pada bab ini berisi kesimpulan dari semua uraian-uraian pada bab-bab sebelumnya dan juga berisi saran-saran yang diharapkan berguna dalam penerapan perangkat lunak deteksi kemiripan dokumen bahasa Indonesia di penelitian selanjutnya.

1.8 Kesimpulan

Penelitian ini berfokus pada deteksi kemiripan dokumen berdasarkan kemiripan penulisan tekstual pada dokumen bahasa Indonesia, serta untuk membandingkan dua algoritma untuk melihat algoritma mana yang lebih baik dalam mendeteksi kemiripan dokumen. Algoritma yang digunakan adalah

algoritma *Jaro-Winkler Distance* dan *Levenshtein Distance*. Dokumen yang digunakan adalah dokumen bahasa Indonesia dengan ekstensi *.pdf* dan *.txt*. Hasil akhir yang diharapkan dari penelitian ini adalah mengetahui perbedaan persentase kemiripan dan waktu proses yang dihasilkan oleh kedua algoritma tersebut.

DAFTAR PUSTAKA

- Adiwidya, B. M. D. (2009). Algoritma Levenshtein Dalam Pendekatan Approximate.
- Adriani, M., Asian, J., Nazief, B., Tahaghoghi, S. M. M., & Williams, H. E. (2007). Stemming Indonesian: A confix-stripping approach. *ACM Transactions on Asian Language Information Processing*, 6(4).
- Agusta, L. (2009). Perbandingan Algoritma Stemming Porter Dengan Algoritma Nazief & Adriani Untuk Stemming Dokumen Teks Bahasa Indonesia. *Konferensi Nasional Sistem Dan Informatika 2009*, (KNS&I09-036), 196–201.
- Friendly, F. (2019). Jaro-Winkler Distance Improvement for Approximate String Search Using Indexing Data for Multiuser Application. *Journal of Physics: Conference Series*, 1361(1).
- Gomma, W. H., & Fahmy, A. A. (2013). A Survey of Text Similarity Approaches. *International Journal of Computer Applications*, 68(13), 13–18.
- Kurniawan, W., Bijaksana, M. A., & Wahyudi, B. A. (2018). Analisis Pencocokan Nama dengan Nama Arab Terjemahan Bahasa Indonesia Menggunakan metode Levenshtein Distance Pendahuluan Studi Terkait Hadits, 5(3), 7472–7493.
- Kurniawati, A. (2010). Implementasi Algoritma Jaro-Winkler Distance untuk

Membandingkan Kesamaan Dokumen Berbahasa Indonesia. *Proceeding, Seminar Ilmiah Nasional Komputer Dan Sistem Intelijen KOMMIT 2008, Depok, Indonesia.*

Rochmawati, Y., & Kusumaningrum, R. (2016). Studi Perbandingan Algoritma Pencarian String dalam Metode Approximate String Matching untuk Identifikasi Kesalahan Pengetikan Teks. *Jurnal Buana Informatika*, 7(2), 125–134.

Sastroasmoro, S. (2007). Beberapa Catatan tentang Plagiarisme *. *Maj Kedokt Indon, Volum: 57*, 239–244.

Su, Z., Ahn, B. R., Eom, K. Y., Kang, M. K., Kim, J. P., & Kim, M. K. (2008). Plagiarism detection using the Levenshtein distance and Smith-Waterman algorithm. *3rd International Conference on Innovative Computing Information and Control, ICICIC'08*, 0–3.

Tinaliah, T., & Elizabeth, T. (2018). Perbandingan Hasil Deteksi Plagiarisme Dokumen dengan Metode Jaro-Winkler Distance dan Metode Latent Semantic Analysis. *Jurnal Teknologi Dan Sistem Komputer*, 6(1), 7.

Yulianto, M. M., Arifudin, R., & Alamsyah, A. (2018). Autocomplete and Spell Checking Levenshtein Distance Algorithm To Getting Text Suggest Error Data Searching In Library. *Scientific Journal of Informatics*, 5(1), 75.